

Combining Statistical Monitoring and Predictable Recovery for Self-Management

Armando Fox
Computer Science
Department
Stanford University

fox@cs.stanford.edu

Emre Kıcıman
Computer Science
Department
Stanford University

emrek@cs.stanford.edu

David Patterson
Computer Science Division
University of California,
Berkeley

patterson@cs.berkeley.edu

ABSTRACT

Complex distributed Internet services form the basis not only of e-commerce but increasingly of mission-critical network-based applications. What is new is that the workload and internal architecture of three-tier enterprise applications presents the opportunity for a new approach to keeping them running in the face of many common recoverable failures. The core of the approach is anomaly detection and localization based on statistical machine learning techniques. Unlike previous approaches, we propose anomaly detection and pattern mining not only for operational statistics such as mean response time, but also for *structural* behaviors of the system—what parts of the system, in what combinations, are being exercised in response to different kinds of external stimuli. In addition, rather than building baseline models *a priori*, we extract them by observing the behavior of the system over a short period of time during normal operation. We explain the necessary underlying assumptions and why they can be realized by systems research, report on some early successes using the approach, describe benefits of the approach that make it competitive as a path toward self-managing systems, and outline some research challenges. Our hope is that this approach will enable “new science” in the design of self-managing systems by allowing the rapid and widespread application of statistical learning theory techniques (SLT) to problems of system dependability.

1. RECOVERY AS RAPID ADAPTATION

A “five nines” availability service (99.999% uptime) can be down only five minutes a year. Putting a human in the critical path to recovery would expend that entire budget on a single incident, hence the increasing interest in self-managing or so-called “autonomic” systems. Although there is extensive literature on statistics-based change point detection [2], some kinds of partial failures, or “brown-outs” in which only part of a service malfunctions, cannot be easily detected by such techniques. For example, one of the au-

thors experienced a bug such that after clicking to purchase a flight for April, a later visit to the “flight details” page showed the wrong flight date (in October) and no flight itinerary details at all. If the operational statistics such as response time for delivering this page are within normal thresholds, performance monitoring would not find this problem. Indeed, such “glitches” accounted for up to 64% of user-visible failures in the 40 top-performing Web sites according to an industry group survey [4], and our own anecdotal experience with one large service provider reveals that 75% of their total recovery time for application-level failures is spent detecting them, and another 18% diagnosing them [9].

We believe a promising direction is to start thinking not in terms of normal operation vs. recovery, but in terms of constant and rapid adaptation to external conditions, including sudden workload changes, inevitable hardware and software failures, and human operator and programmer errors. In particular, we propose the broad application of techniques from statistical learning theory (SLT)—automatic classification, novelty/anomaly detection, data clustering, etc.—to observe and track *structural* behaviors of the system, and to detect potential problems such as the example above.

2. APPROACH AND ASSUMPTIONS

We assume typical request-reply based Internet services, with separate session state [18] used to synthesize more complex interactions from a sequence of otherwise stateless request-reply pairs. Past approaches to statistical monitoring of such services have primarily relied on *a priori* construction of a system model for fault detection and analysis; this construction is tedious and error-prone, and will likely remain so as our services continue to evolve in the direction of heterogeneous systems of black boxes, with subsystems such as Web servers, application logic servers, and databases being supplied by different vendors and evolving independently. We propose instead to build and periodically update the baseline model by observing the system’s own “normal” behavior. The approach can be summarized as follows:

1. Ensure the system is in a state in which it is mostly doing the right thing most of the time, according to simple and well-understood external indicators.
2. Collect observations about the system’s behavior during this time to build one or more baseline models

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

WOSS’04 Oct. 31–Nov. 1, 2004 Newport Beach, CA, USA
Copyright 2004 ACM 1-58113-989-6/04/0010 ...\$5.00.

of behavior. These models may capture either time-series behaviors of particular parameters or structural behaviors of the system.

3. If “anomalous” behaviors relative to any of these models are observed, automatically trigger simple corrective actions. If repeated simple corrective actions do not cause the anomaly to go away, notify a human. Since false positives are a fact of life with statistical approaches, we also need a strategy for quantifying and dealing with the cost of acting on false positives.
4. Periodically, go back to step 2, to update the model.

Each of steps 1–3 corresponds to an assumption, as follows.

A1. Large number of independent requests. If most users’ interactions with the service are independent of each other (as they usually are for Internet services), and if we assume bugs are the exception rather than the norm, such a workload gives us the basis to make “law of large numbers” arguments supporting the use of statistical techniques to extract the model from the behavior of the system itself.

A2. Ability to rapidly characterize “normal” behavior through observation of relevant features. To use statistical or data-mining techniques, we need a representation of the data observations the model will operate on (“concepts” in the terminology of data mining) and a non-intrusive way to capture those observations. For example, a particular codepath through a componentized application will “touch” some subset of the components, and it is likely that the number of legitimate code paths is much smaller than the number of permutations of components, making anomaly detection appealing. Further, in order to build a baseline (“normal”) model rapidly, the behaviors of interest must be exercised often enough and stable for long enough that the time needed to train a model is significantly shorter than the time over which the model is expected to be useful. Busy Internet services serve a large number of users per unit time, so that large fractions of the service’s functionality are exercised in a relatively short period, providing hope that the model can be created and maintained online while the system is running.

A3. Simple and predictable control points. If the model’s predictions and analyses are to be used to effect service repair when an anomaly indicating a potential failure is detected, there must be a safe, predictable, and relatively non-disruptive way to do so. *Safe* means that correct application semantics are not jeopardized by actuating the control point. *Predictable* means that the cost of actuating the control point must be well known. *Non-disruptive* means that the result of activating a control point will be no worse than a minor and temporary effect on performance. These properties are particularly important when statistical techniques are used because those techniques will inevitably generate false positives. If we know that the only effect of acting on a false positive is a temporary and small decrease in performance, we can quantify the cost of “blindly” acting on false positives; this enhances the appeal of automated statistical techniques, since many techniques’ sensitivity can be tuned to trade off false positive rates vs. false negative (miss) rates.

We now turn to how these assumptions might be satisfied in a real service. Note that A1 is trivially true for the

services in question, whereas A2 and A3 lead to some interesting systems research.

3. OBSERVATION AND CONTROL POINTS

Since modifying every existing application to add observation and control points is cumbersome and unlikely, we limit our attention initially to *framework-intensive* applications¹—those whose total code consists mostly of middleware (e.g. J2EE runtime services, libraries, etc.) with a smaller amount of application logic (though even simple applications typically contain 10K to 100K lines of such logic). By modifying the middleware, we can provide application-generic observation points without any extra work for application programmers. For example, we modified the source code of the JBoss open-source application server to collect and report code-path observations [8].

It is more difficult to add application-generic control points that are predictable, safe and non-disruptive. Crashing and rebooting a machine is certainly predictable, since crashing relies only on a simple external mechanism (the power switch), but it may be unsafe or disruptive or both, unless the application is known to be crash-only [5]. We have successfully done this by modifying a J2EE application server to support the “microrebooting” of individual application components [6] and have begun combining this ability with application-generic statistical learning based failure detection [7]. Another alternative is coercing any failure to a machine-level crash and turning the failure into slight additional latency using a combination of overprovisioning and fast failover, as is commonly done for stateless Web servers [3] and more recently for specialized state storage subsystems [17, 13].

We give examples using microrebooting since we have some early results in that area, but other methods of microrecovery are being independently explored, such as the unloading and reloading of kernel device drivers to recover from transient driver-related failures [21].

4. SLT AND DEPENDABILITY

Having briefly addressed some important systems-building issues (to which we return shortly in the context of some concrete examples), we now discuss the core of our detection and diagnosis strategy. Statistical learning theory (SLT) provides a framework for the design of algorithms for classification, prediction, feature selection, clustering, sequential decision-making, novelty detection, trend analysis, and diagnosis. Its techniques are already being used in bioinformatics, information retrieval, spam filtering and intrusion detection. We propose a software architecture for integrating SLT pervasively into the computing infrastructure, as a tool for evaluating which SLT techniques are useful at detecting which kinds of problems. For concreteness, we describe two simple examples: one based on time-series models and another based on structural models.

4.1 Time Series Models

Time-series models capture repeatable patterns in a service’s temporal behavior. For example, the memory used by a server-like process typically grows until garbage collection occurs, then falls abruptly. We do not know the period of

¹By framework we refer to a componentized middleware such as J2EE or CORBA.

this pattern, or indeed whether it is periodic; but we would expect that multiple servers running the same logic under reasonable load balancing should behave about the same—the relative frequencies of garbage-collection events at various timescales should be comparable across all the replicas. We successfully used this method to detect anomalies in replicas of SSM, our session state management subsystem [18]. Each replica reports the values of several resource-usage and forward-progress metrics once per second, and these time series are fed to the Tarzan algorithm [14], which discretizes the samples to obtain binary strings and counts the relative frequencies of all substrings within these strings. Normally, these relative frequencies are about the same across all replicas, even if the garbage-collection cycles are out of phase or their periods vary². If the relative frequencies of more than 2/3 of these metrics on some replica differ from those of the other replicas, that replica is immediately re-booted. Since SSM’s replication-based design has some inherent overprovisioning and each replica is optimized for fast reboot, rebooting is safe, predictable and non-disruptive. Indeed, SSM has no concept of “recovery” vs “normal” behavior; since periodic reboots are normal and incur little performance cost, the system is “always recovering” by adapting to changing external conditions through a simple composition of mechanisms.

We point out that we chose a collection of metrics to monitor that, according to our understanding of the application, should have some correlation with successful forward progress. In practice, we found that either almost none of the metrics exhibit anomalies, or most of them do, hence the 2/3 threshold. Although a less naive choice of algorithm or more careful selection of metrics might be more efficient, we are encouraged by the baseline results we obtained even without such optimizations.

4.2 Structural Models

Structural models capture control-flow behavior of an application, rather than temporal behavior. Whereas we applied a simple time-series method above to an application we built ourselves, we have found structural models useful for characterizing the behavior of an application whose structure is *not* well understood in advance. One example of a structural model is a *path*—the inter-component dynamic call tree resulting from a single request-reply interaction. We modified an open-source J2EE application server, JBoss, to dynamically collect such call trees for all incoming requests; these are then treated as parse trees generated by a probabilistic context-free grammar (PCFG) [19]. Later on, when a path is seen that corresponds to a low-probability parse tree, the corresponding user request is flagged as anomalous. The scoring function, described in [15], compares the expected probability of a given inter-component call being part of a path with the observed frequency of that call within that path at runtime. Overall, this approach detects over 90% of various injected faults. More interestingly, however, figure 1 shows that increasing the algorithm’s sensitivity parameter α reduces the number of false positives (better precision), but increases the false negative rate (i.e., decreases the overall detection rate from an average of 68% to an average of 34%) and increases the number of experiments in which no failures are detected at all despite fault injection

²Classical time-series methods are less effective when the signal period varies.

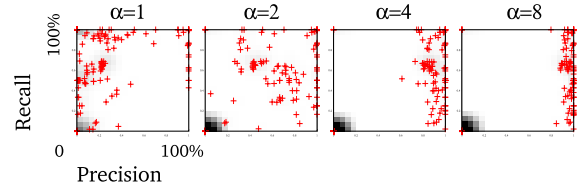


Figure 1: Detection rate (inverse of false negative rate) vs. precision (inverse of false positive rate) for PCFG-based path-shape analysis of PetStore 1.3 running on our modified JBoss server. Each point represents one experimental run, with shaded tiles representing clustering of points too densely to depict each one individually. As sensitivity is increased, false positives are reduced, but overall detection rate goes down (failures go undetected in a larger number of experiments, as shown by the clustering of points in the lower-left corner), and even in experiments in which failures are detected, the average detection rate goes down.

(i.e. both recall and precision are zero). Hence, in order to minimize false negatives (maximize recall), we must be willing to accept nontrivial false positive rates.

The key to our approach is that by making our control actions non-disruptive and safe, we can in fact tolerate significant false positive rates. In this case, we respond by selectively “microrebooting” the suspected-faulty EJB’s without causing unavailability of the entire application. Although this work is still in progress, we have demonstrated that EJB microreboots are predictable and non-disruptive, allowing false positive rates as high as 97% while still incurring less downtime than a corresponding full reboot [6]. Microreboots in J2EE are safe because J2EE constrains application structure in a way that makes most persistent state management explicit, allowing us to externalize the session state into SSM so that it survives microreboots of the EJB’s.

5. DISCUSSION

The combination of low-cost recovery with statistical anomaly detection raises several interesting additional discussion points; we focus on three that we believe arise from the novelty of the approach.

False positives don’t matter. All statistical techniques are prone to some level of false positives, and a typical trade-off in anomaly detection algorithms involves the detection rate (of true failures) vs. the false positive rate (of events misidentified as anomalous). Historically, minimizing the false positive rate has been a major goal of algorithm designers. However, by making the cost of acting on a false positive sufficiently low, not only can we tolerate false positives, but we can potentially combine overlapping analysis techniques to improve overall coverage and detection.

Correlation, causation, and diagnosis. We avoid describing our approach as “diagnosis,” both because diagnosis requires application-specific information or human programmer effort and because techniques such as rebooting illustrate that diagnosis is not always a prerequisite to repair. Rather, we identify events or behaviors that are anomalous with respect to failure-free operation, narrow down to the specific components or subsystems correlated

with the anomaly (we might optimistically call this “localization”), and invoke low-cost techniques that have a good track record of repairing certain classes of problems whose true cause may be unknown. One challenge arising from this approach is systematic misdiagnosis (and the resulting corrective action) of rare-but-legitimate events, such as an infrequently-traversed but valid codepath. Even low-cost recovery, if invoked often enough, begins to interfere with normal operation. We will address this challenge systematically but we believe part of the solution involves combining a number of detection and monitoring methods rather than relying exclusively on the techniques proposed here.

Application structure. For correctness, we exploited the fact that today’s Internet applications nicely separate process recovery from data recovery; in effect, the session state manipulated between stateless HTTP requests amounts to a “microcheckpoint” of the application, so that component-level recovery occurring between microcheckpoints is correctness-preserving. Nonetheless, the approach is more generally applicable and is synergistic with projects such as ARMOR [24], which allows retrofitting existing applications with a “microcheckpoint” facility.

6. RESEARCH CHALLENGES

We have focused on recasting “recovery” as a kind of rapid adaptation, but a similar argument applies for other online operations such as resource management. For example, online repartitioning of a cluster-based hash table [13] can be achieved by taking one replica offline (which looks like a failure and does not affect correctness), cloning it, and bringing both copies back online. The resulting stale data is automatically repaired by normal-case mechanisms, hence no new machinery is required to implement incremental scaling as an online operation. This is a fine-grained analog to similar techniques used in very large scale systems, e.g. a server farm designed to handle partial failure can be migrated in two parts by changing the DNS pointer while half the farm is still operating at each physical site [3].

Most existing implementations of SLT algorithms are offline; our proposal may motivate SLT developers to focus on online and distributed algorithms. The above experiments show that even an unoptimized offline implementation of PCFG analysis can process thousands of paths in a few seconds. This in turn motivates the need for generic data collection and management architectures for statistically-monitored systems: even a simple (11K lines of code) application we instrumented produces up to 40 observations per user request, with 1,000 to 10,000 requests per second being representative of Internet services. Scalable abstractions for sliding data windows, sampling, fusion of results from different SLT models, etc. will have to be provided, as well as easy ways to create observation and control points without requiring intrusive modifications to every application.

Finally, although we have discussed applying SLT approaches primarily at the application level, we note that the needed infrastructure is largely in place for applying it at all levels of functionality all the way down to the hardware. A legacy of the Active Networking research agenda [22] is a new generation of user-programmable network devices for storage virtualization, server load balancing, and traffic management, which provide some of the observation and control points needed by our approach and allow us to make “law of large numbers” arguments required by assumption

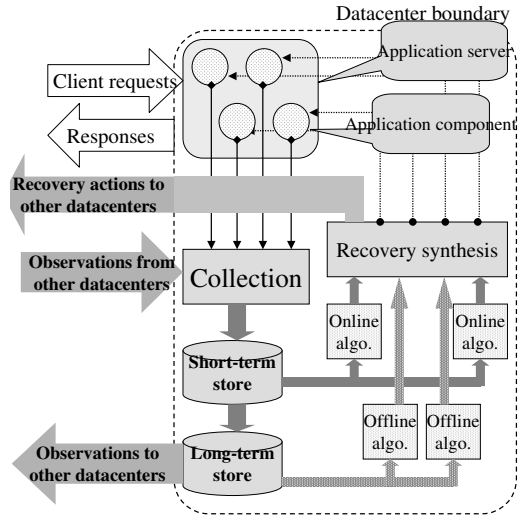


Figure 2: High-level architecture of an intra-datacenter adaptive system. We have an early prototype of the application server, have acquired some initial experience with online detection algorithms, and have explored simple recovery synthesis based on statistical localization and microrebooting.

A1. Figure 2 shows a block-diagram architecture (parts of which we are already prototyping) for distributed network applications that exploit SLT-based monitoring at multiple levels.

7. RELATED WORK

Anomaly detection and classification have been used to infer errors in systems code [10], debug Windows Registry problems [23], detect possible violation of runtime variable assignment invariants [12], and discover source code bugs by distributed assertion sampling [16]. The latter is particularly illustrative of SLT’s ability to mine large quantities of observations for interesting patterns that can be directly related to dependability. System parameter tuning and automatic resource provisioning have also been tackled using PCFG-based approaches [1] and closed-loop control theory [20], although such approaches generally cannot detect functional or structural deviations in system behavior unless they manifest as performance anomalies. Finally, the Recovery-Oriented Computing project [11] has argued that fast recovery is good for its own sake, but in the context of SLT, fast recovery is *essential* because it gives us an inexpensive way to deal with false positives. As such, ROC is a key enabler for this approach, and we build on its success.

8. CONCLUSION

Our ability to design and deploy large complex systems has outpaced our ability to deterministically predict their behavior except at the coarsest grain. We believe statistical

approaches, which can find patterns and detect deviations in data whose semantics are initially unknown, will be a powerful tool not only for monitoring and online adaptation of these systems but for helping us better understand their structure and behavior. In particular, by fitting applications with low-cost “microrecovery” actions, we can tolerate nontrivial false positive rates in statistical detection algorithms, allowing us to improve their sensitivity and coverage to the point where they can detect failures that would go unnoticed by existing monitoring techniques. We believe our encouraging initial results with J2EE applications and simple well-known machine learning algorithms invite deeper exploration of the approach. A generic platform for pervasive integration of SLT methods, themselves the subject of broad and vigorous research, would hasten the adoption of SLT into dependable systems, which we believe would in turn provide a new scientific foundation for the construction of self-managing systems.

Acknowledgments

The ideas in this paper have benefited from advice and discussion with George Candea, Timothy Chou, Moises Goldszmidt, Joseph L. Hellerstein, Ben Ling, Matthew Merzbacher, and Chris Overton, and many others.

9. ADDITIONAL AUTHORS

Additional authors: Randy Katz, Michael Jordan, University of California, Berkeley (email: {randy, jordan}@cs.berkeley.edu)

10. REFERENCES

- [1] Paul Barham, Rebecca Isaacs, Richard Mortier, and Dushyanth Narayanan. Magpie: real-time modelling and performance-aware systems. In *Proc. 9th Workshop on Hot Topics in Operating Systems*, Lihue, HI, June 2003.
- [2] Michele Basseville and Igor V. Nikiforov. *Detection of Abrupt Changes—Theory and Application*. Prentice-Hall Inc., Englewood Cliffs, NJ, 1993.
- [3] Eric Brewer. Lessons from giant-scale services. *IEEE Internet Computing*, 5(4):46–55, July 2001.
- [4] Business Internet Group. The black Friday report on Web application integrity. San Francisco, CA, 2003.
- [5] George Candea and Armando Fox. Crash-only software. In *Proc. 9th Workshop on Hot Topics in Operating Systems*, Lihue, HI, June 2003.
- [6] George Candea, Shinichi Kawamoto, Yuichi Fujiki, Greg Friedman, and Armando Fox. A microbootable system – design, implementation, and evaluation. In *Proc. 6th USENIX Symposium on Operating Systems Design and Implementation*, San Francisco, December 2004.
- [7] George Candea, Shinichi Kawamoto, Emre Kiciman, and Armando Fox. Autonomous recovery in componentized internet applications. *Cluster Computing Journal*, To appear, 2005.
- [8] George Candea, Pedram Keyani, Emre Kiciman, Steve Zhang, and Armando Fox. JAGR: An autonomous self-recovering application server. In *Proc. 5th International Workshop on Active Middleware Services*, Seattle, WA, June 2003.
- [9] Mike Chen, Emre Kiciman, Anthony Accardi, Armando Fox, and Eric Brewer. Using runtime paths for macro analysis. In *Proc. 9th Workshop on Hot Topics in Operating Systems*, Lihue, HI, 2003.
- [10] Dawson Engler, David Yu Chen, Seth Hallett, Andy Chou, and Benjamin Chelf. Bugs as deviant behavior: A general approach to inferring errors in systems code. In *Proc. 18th ACM Symposium on Operating Systems Principles*, pages 57–72, Lake Louise, Canada, Oct 2001.
- [11] David A. Patterson et al. Recovery-Oriented Computing: motivation, definition, techniques, and case studies. Technical Report CSD-02-1175, University of California at Berkeley, 2002.
- [12] Sudheendra Hangal and Monica Lam. Tracking down software bugs using automatic anomaly detection. In *Proceedings of the International Conference on Software Engineering*, May 2002.
- [13] Andrew C. Huang and Armando Fox. Cheap recovery: A key to self-managing state. Submitted for publication.
- [14] E. Keogh, S. Lonardi, and W Chiu. Finding surprising patterns in a time series database in linear time and space. In *In proc. of the 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 550–556, Edmonton, Alberta, Canada, Jul 2002.
- [15] Emre Kiciman and Armando Fox. Detecting application-level failures in component-based internet services. Submitted for publication, September 2004.
- [16] Ben Liblit, Alex Aiken, Alice X. Zheng, and Michael I. Jordan. Bug isolation via remote program sampling. In *Proceedings of the ACM SIGPLAN 2003 Conference on Programming Language Design and Implementation*, San Diego, California, June 9–11 2003.
- [17] Benjamin Ling, Emre Kiciman, and Armando Fox. Session state: Beyond soft state. In *Proc. 1st Symposium on Networked Systems Design and Implementation*, San Francisco, CA, March 2004.
- [18] Benjamin C. Ling, Emre Kiciman, and Armando Fox. Session state: Beyond soft state. In *Proc. 1st USENIX/ACM Symposium on Networked Systems Design and Implementation*, San Francisco, CA, March 2004.
- [19] Christopher D. Manning and Hinrich Shutze. *Foundations of Statistical Natural Language Processing*. The MIT Press, Cambridge, MA, 1999.
- [20] S Parekh, N Gandhi, JL Hellerstein, D Tilbury, TS Jayram, and J Bigus. Using control theory to achieve service level objectives in performance management. *Real Time Systems Journal*, 23(1–2), 2002.
- [21] Michael Swift, Brian N. Bershad, and Henry M. Levy. Improving the reliability of commodity operating systems. In *Proc. 19th ACM Symposium on Operating Systems Principles*, 2003.
- [22] David L. Tennenhouse and David J. Wetherall. Towards an active network architecture. In *ACM SIGCOMM '96 (Computer Communications Review)*. ACM, 1996.
- [23] Yi-Min Wang, Chad Verbowski, and Daniel R. Simon. Persistent-state checkpoint comparison for troubleshooting configuration failures. In *Proc. International Conference on Dependable Systems and Networks*, San Francisco, CA, June 2003.
- [24] K. Whisnant, S. Bagchi, B. Srinivasan, Z. Kalbarczyk, and R.K. Iyer. Incorporating reconfigurability, error detection, and recovery into the Chameleon ARMOR architecture. Technical Report CRHC-98-13, University of Illinois at Urbana-Champaign, 1998.